# Cyber Tools for Large Collaborations

James D. Myers

jimmyers@ncsa.uiuc.edu
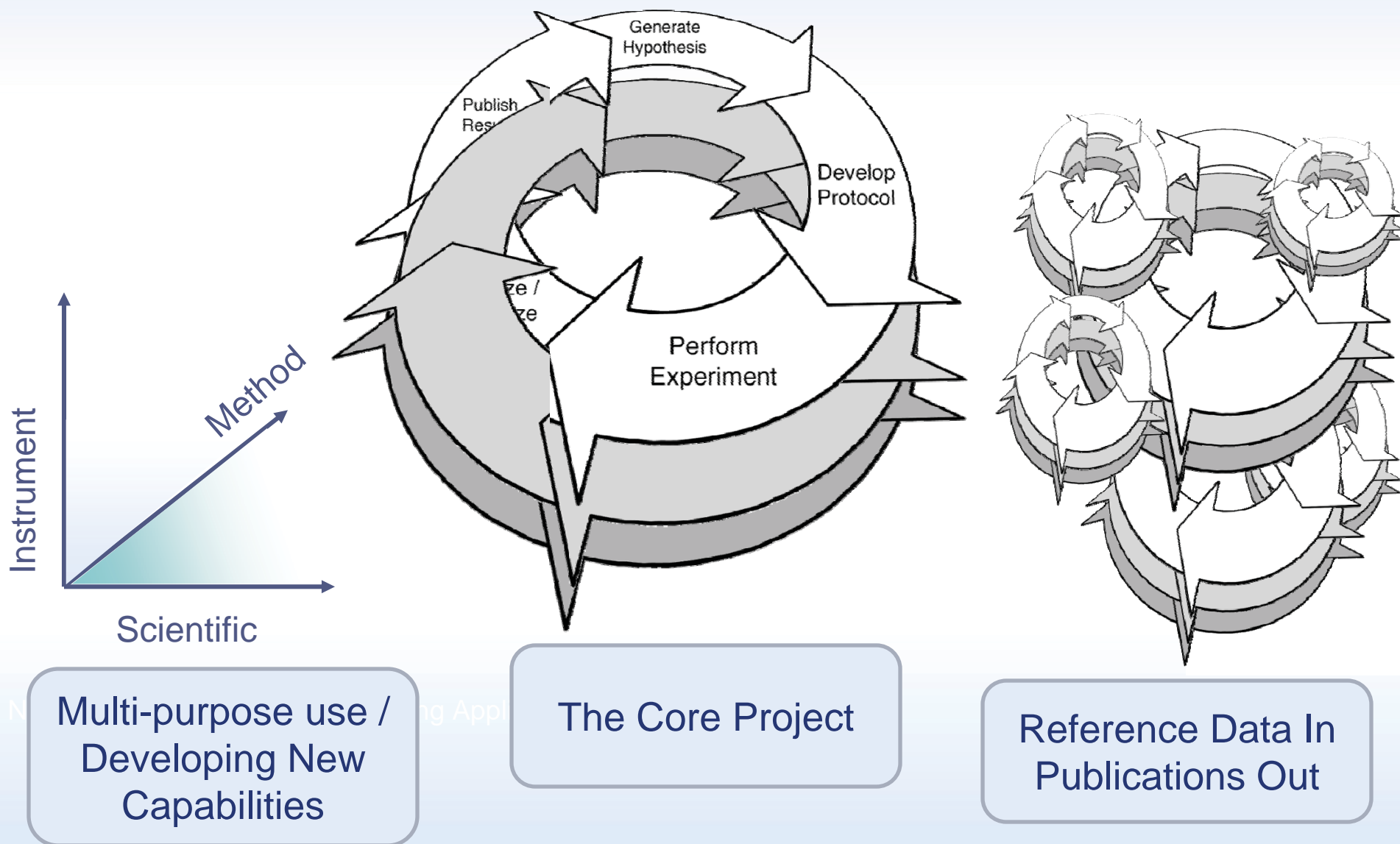
National Center for Supercomputing Applications
University of Illinois at Urbana-Champaign
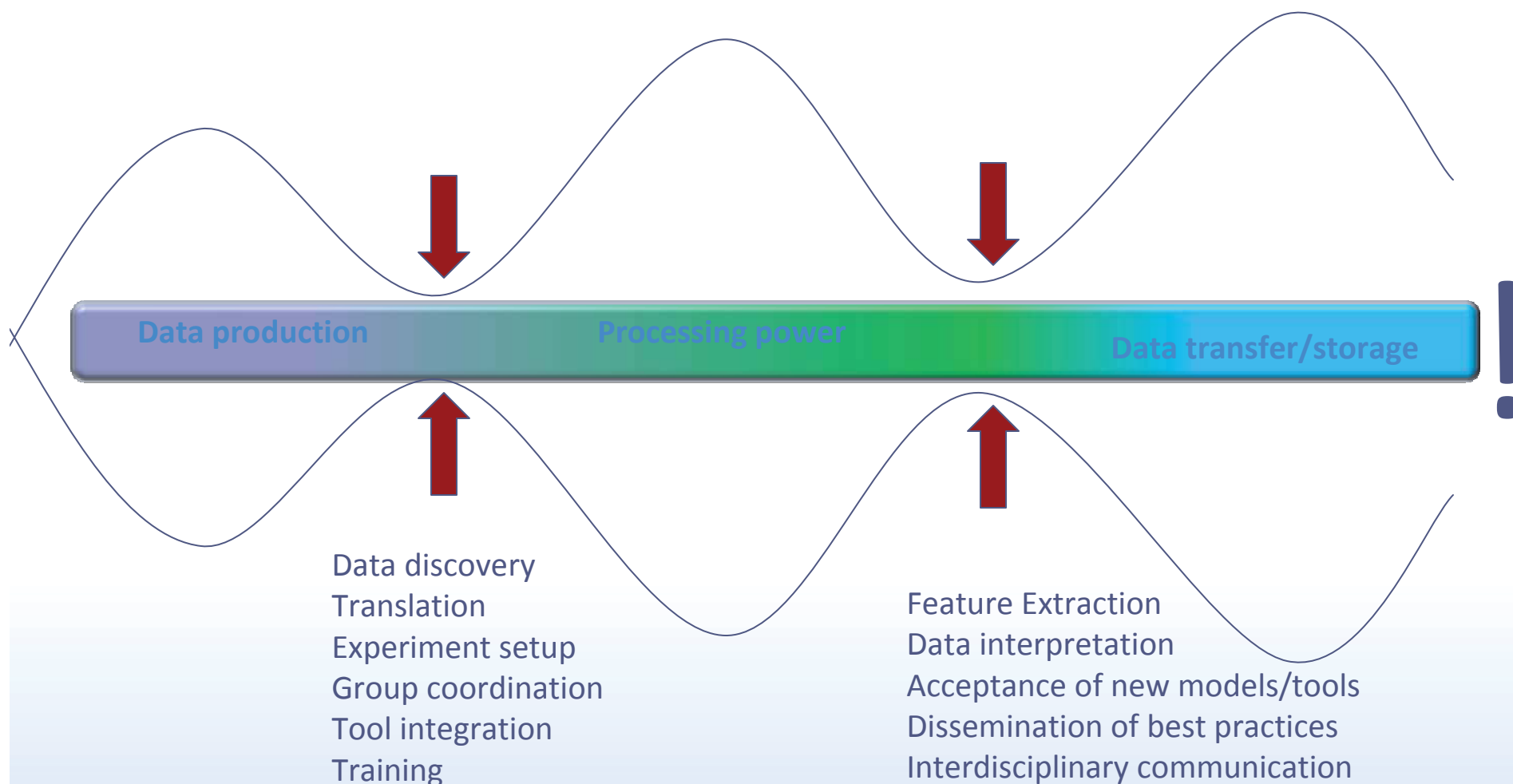
## Outline

- What do we need Cyber Tools to do?
- What forms of Scalability do we need?
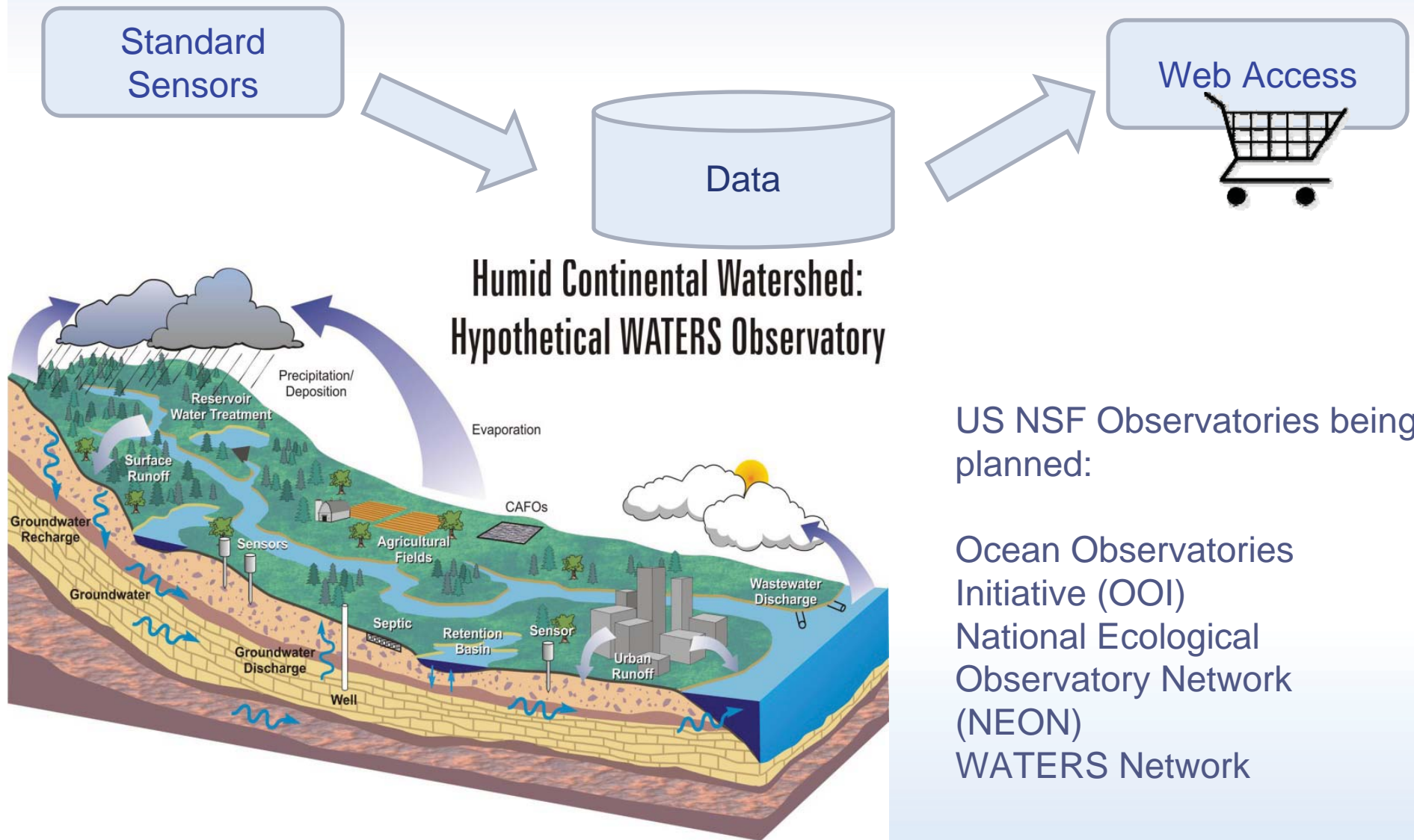- Design 'patterns' for building scalable tools
- Some examples…

# The Research Process



Instrument
Method
Scientific

Generate Hypothesis
Develop Protocol
Perform Experiment
Publish Res...
...ze / ...ze

Multi-purpose use / Developing New Capabilities

The Core Project

Reference Data In Publications Out

# Cyberenvironments Recognize 'Amdahl's Law' for Scientific Progress



Data production   Processing power   Data transfer/storage

!

Data discovery
Translation
Experiment setup
Group coordination
Tool integration
Training

Feature Extraction
Data interpretation
Acceptance of new models/tools
Dissemination of best practices
Interdisciplinary communication

# Environmental Observatories: Dissemination Model



Standard Sensors → Data → Web Access

Humid Continental Watershed: Hypothetical WATERS Observatory

Precipitation/Deposition
Reservoir Water Treatment
Surface Runoff
Groundwater Recharge
Groundwater
Groundwater Discharge
Well
Sensors
Septic
Retention Basin
Agricultural Fields
Sensor
Evaporation
CAFOs
Urban Runoff
Wastewater Discharge

US NSF Observatories being planned:

Ocean Observatories Initiative (OOI)
National Ecological Observatory Network (NEON)
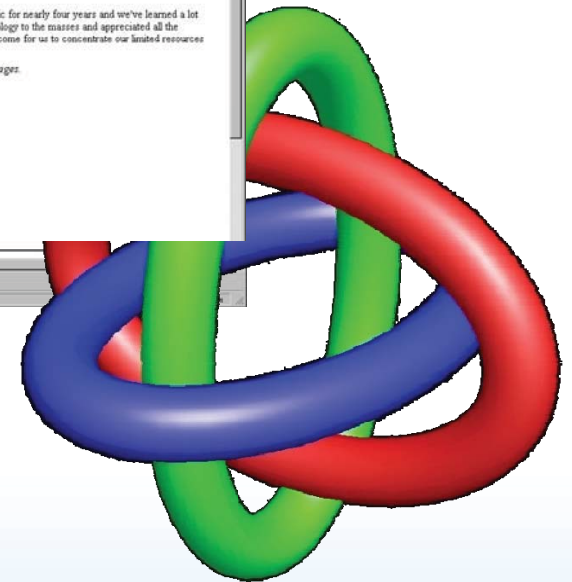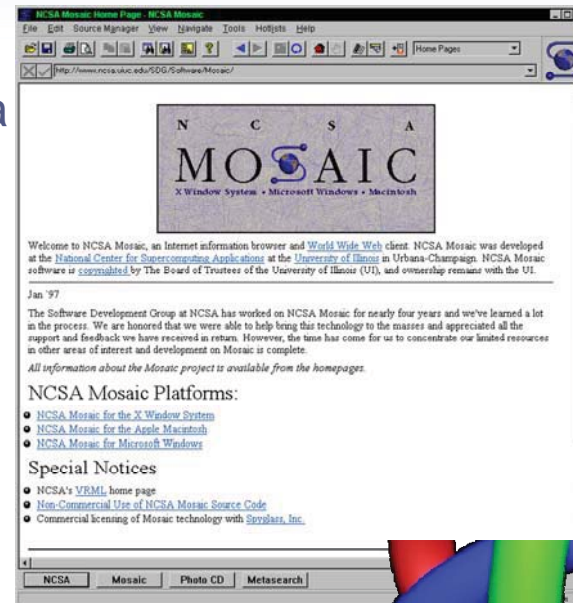WATERS Network

NCSA

# Virtual Observatories

You and your colleagues are there:

...your view
...your data
...your models
...your notes, papers

DNA Micro-Arrays

Chem-Lab on a chip

Minirhizotron Array

Micro-weather Stations

Vapor Detector

E-nose

Smart Sensor Web

Humi

NCSA
CYBERCOLLABORATORY

Welcome, James D Myers!

My Account
SIGN OUT

My Menu
My Posts
My Documents

Automated E-ton

Precipitation/ Deposition

Reservoir Water Treatment

Surface Runoff

Groundwater Recharge

Sensors

Groundwater

Agricul Field

Septic

Basin

Well

Groundwater Discharge

Wastewater Discharge

Urban Runoff

27.82
27.8
27.78
27.76
27.74
27.72
27.7

-97.28   -97.26   -97.24   -97.22   -97.2   -97.18   -97.16   -97.14

0   0.1   0.2   0.3   0.4   0.5   0.6   0.7   0.8   0.9

NCSA

# *Is the World Wide Web a way to share web pages?*

- Mosaic
  - By early 1990s, the internet had a wealth of resources, but they were inaccessible to most scientists
  - *Individual publishing*
  - *Browsing versus retrieving*
  - *See "Web 2.0 ... The Machine is Us/ing Us"*

- Cyberenvironments
  - By the early 2000's, the internet and grid had a wealth of interactive resources, but they were inaccessible to most scientists
  - *Individual information models*
  - *Fusion versus gathering*



See "The Machine is Us/ing Us"! Michael Wesch

NCSA

# Beyond Data 'Take-Out':

- **Enable synthesis of multiple types of data?**
- **Provide useful statistical and visual summaries?**
- **Combine observational and modeled data?**
- **Integrate derived products from within and across large heterogeneous edge-less communities?**
- Capture processes for reuse?
- Convey expertise as well as raw resources?
- Enable individuals to create derived data and capabilities that are 'first class citizens'?
- Support rapid dissemination and evolution of preliminary results?
- Enable problem-focused Collaboration?
- Support long-term curation and preservation by third parties?

# Can we build it?

- There are certainly research issues in providing such rich capabilities cost-effectively…
- Can we architect to allow 'innovation at the edges'
    - What do we need to standardize?
    - What do we need to decouple?

Lifecycle cost/benefit analysis

# people

level of capability

Architecting and managing for CI that is "grown, not built"*

*Ixchel Faniel, U. Michigan

**NCSA**

# Relevant Design Patterns for Scalable Not Scaled Cyber Tools

- Abstract interfaces that separate how from what:
  - Authentication ala JAAS, PAM (via callbacks)
  - Content management (via metadata/typed blob abstraction)
  - Global identifiers and declarative semantics (via semantic web)
  - Process abstraction (via workflow and provenance services)
  - Interface integration (via plug-ins, widget, portlets, mash-ups)
  - Event integration (e.g. via enterprise service bus)
  - Virtualization (e.g. services, virtual machines, Grid)

- Standardization occurs via social processes rather than technology lock-in – think TCP-IP, HTTP, XML
- Community-centric projects can help standardize and coordinate (e.g. FEON, Provenance Challenge)

NCSA

# MAEViz: Consequence-Based Risk Management for Seismic Events



- **Engineering View of MAE Center Research**
- Physical through Socio-economic Analysis
- A "Cyberinfrastructure Aware" Application

Decision Support

Damage Prediction

Fragility Models

Inventory Selection

Hazard Definition

## "I have sensitive data I won't distribute"

→ **Network Aware**



- WebDAV, JCR, RDF, SAM, Tupelo

**"Understanding the Scientific Basis of Decisions is Critical"**
**"My calculations are getting too large for my desktop"**

→ **Process Aware**



Process
Capture

Discover

Execute

Report

- Workflow, Provenance, OPM, RDF

**NCSA**

# "Developing a Scenario requires a wide range of expertise"
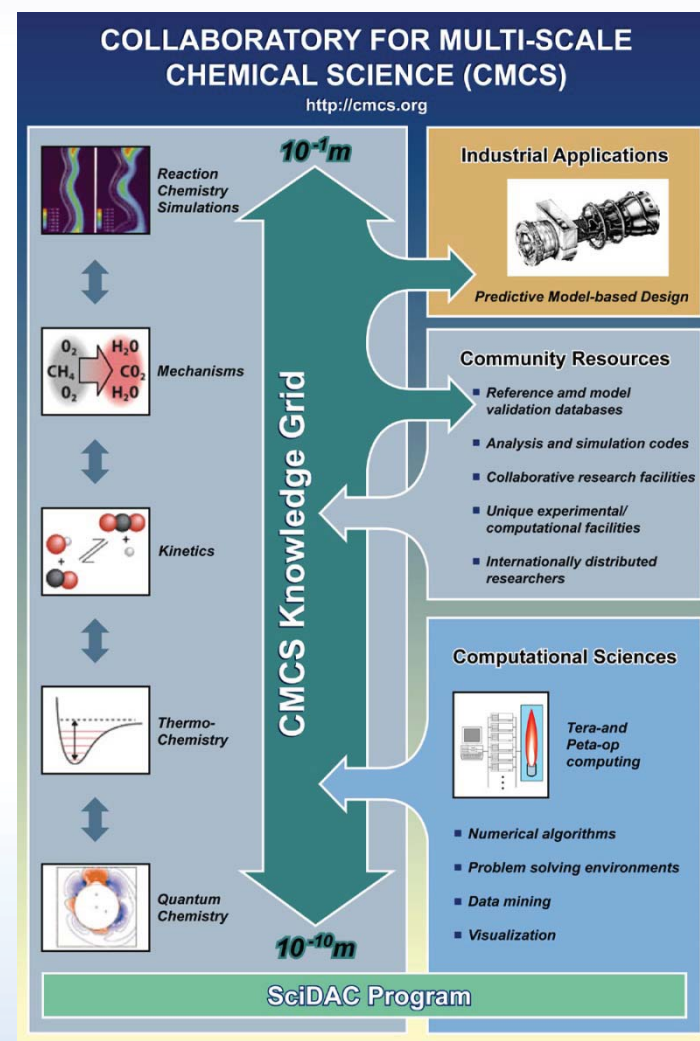
## → Group Aware



Plan, Coordinate,
Share, Compare

SSO

Wiki
Task List
Chat
Document Repository
Scenario Repository
Training Materials

- Collaboratory, Portal, …

NCSA

# "My results could impact how we prepare for the next event"

→ **Dynamic**

New Third-Party
Analyses

Compare, Contrast,
Validate

Auto-update

MAEviz

| Workflow | Data |

Eclipse RCP

• Plug-ins, Provenance, Environment

Plug-in Framework

NCSA

# *Collaboratory for Multiscale Chemical Science*
## *( http://cmcs.org )*

- **A systems-science approach to addressing complex problems**
  - **New knowledge is assimilated from different data, tools, and disciplines at each scale**
  - **Real-time bi-directional information flow**
  - **Multiple applications of the same information**
  - **Evolving scientific models and tools**
- **A cyberenvironment approach:**
  - **General content store with configurable translations**
  - **Publish/subscribe messaging**
  - **Portal, application, and service interfaces**
  - **Multiscale provenance**
  - **+ Standard data/tool/collaboration access**
  - **Community/group coordination, data curation, and model validation**



COLLABORATORY FOR MULTI-SCALE CHEMICAL SCIENCE (CMCS)
http://cmcs.org

# Active ThermoChemical Tables (Ruscic)

- Statistical Analysis of Thermochemical Networks
- Inconsistent inputs
- Simultaneous analysis of uncertainties across network
- 'What If' Analysis
- Potential to Trigger New Calculations

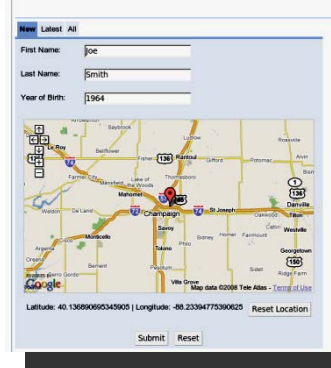# "Range Identification & Optimization Tool" for Mechanism Reduction (Green)



- Mechanism Reduction with Guaranteed Range of Validity
- Web Service @ MIT
- Portlet interface within CMCS
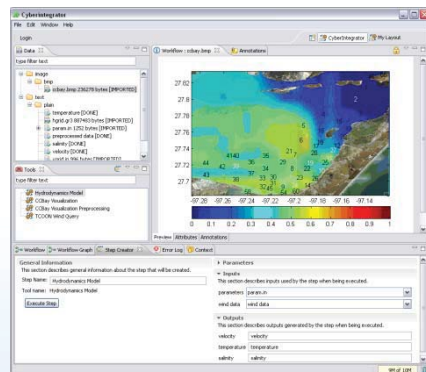- Seamless data transfer
- Asynchronous operation

NCSA

# NCSA's Digital Synthesis Framework (DSF)



**Create**
**Publish**
**Explore**

Historical Data

How does agriculture in the Midwest contribute to fish kills in the Gulf of Mexico?

Where should we send students today to make additional observations?

What are the multi-year trends in hypoxia?

NCSA

# Core DSF Concept

**Web Inputs**

**Workflow Execution Service**

**Visualized Outputs**

Publish

- Parameters
- Input Streams
- Trigger Conditions
- Visualization
- Provenance and
  Annotation Options

**Semantic Content Repository & Provenance Store**

VM Farm /
Compute Cloud

**Desktop Exploration**

NCSA

# Tupelo II: Semantic Content Management

- Web Protocol to
  - Authenticate
  - Get/Set Data
  - Get/Set Metadata
- Flexible Global Identifiers
- Extensions to support specific ontologies (provenance, data streams, GIS, …)

Data Streams

Data Files, Documents

GIS Structures, Images, graphs

Metadata, Provenance

Tupelo Semantic Content Middleware



Local and Distributed Data Sources

# CyberIntegrator:
# Stream-aware Exploratory Workflow



- Identify Inputs
  - (e.g. temperature for the last 24 hours from sensors in a region)
- Link analyses and models
  - Could be Excel, Matlab, or high-end models
- Create/explore visualizations

# Web publication of data and models

- Select Data and Visualizations to Create/Display
- Select Model Inputs (If Any) to Expose
- Publish
  - Register workflow as a service
  - Verify Data/Models Available
  - Generate Input Page and Output Pages as Needed
- Use
  - Web site with dynamic widgets
  - Derived Data in new workflows
  - Widgets embeddable in other pages
  - Annotate, cite, reuse



Desktop Workflow → Hosted Service →

Shared Content Store

NCSA

# Historic Weather Data (e.g. on your Birthday)

# Historic Weather Data

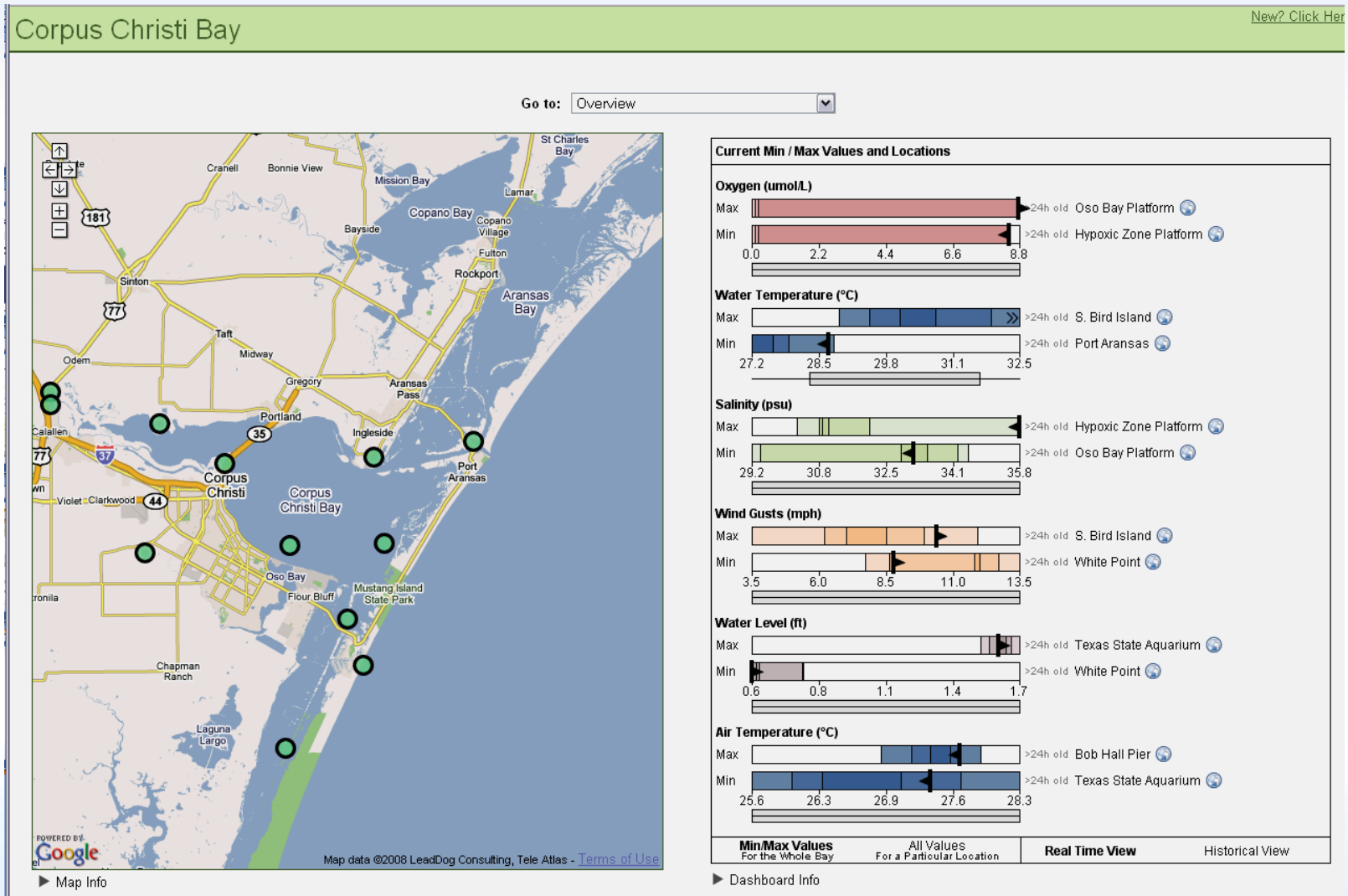# Institute for Genomic Biology and 4-H: Plant Growth Model



Making state-of-the-art plant growth models available for 4-H/student use

Integrating sophisticated modeling into "Seeds and Soils" 4-H activities
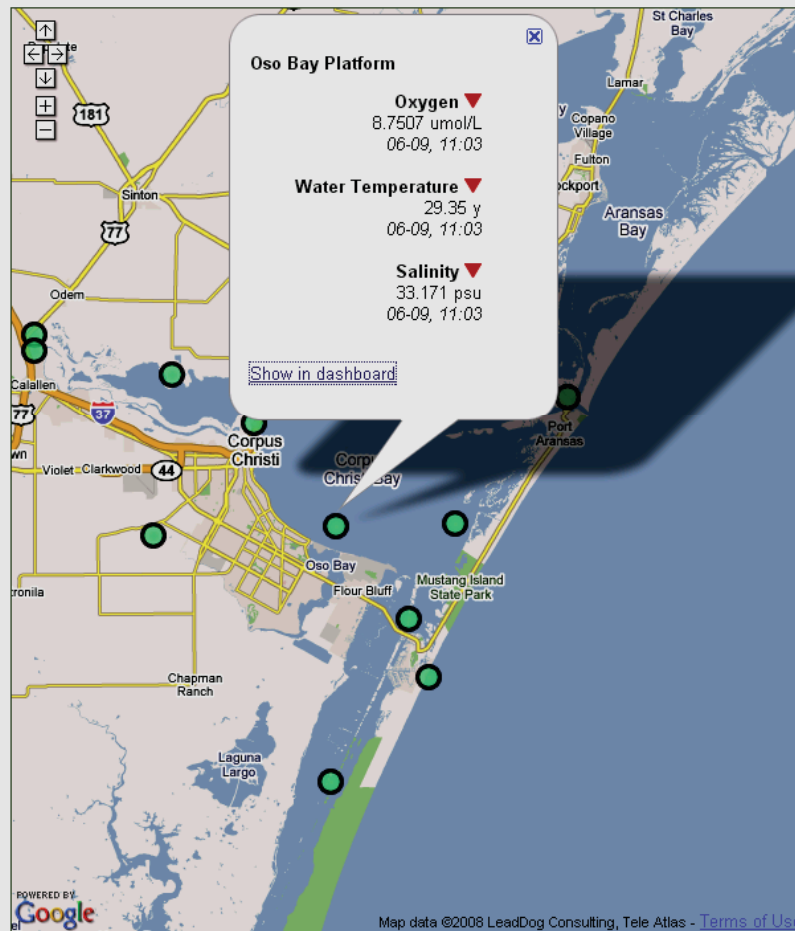
# Corpus Christi Bay Dashboard

# Corpus Christi Bay Dashboard

# Virtual Rain Gauges from Radar Reflectivity

# Impacts

- Mechanism to disseminate your data/model integrated with core Observatory data

- Adds Problem-Specific Capabilities to Observatory and a mechanism for continuing extension and evolution of Observatory services

- Publication allows users to convey expertise

  - (range of validity, 'best' values)

- Data and visualizations can be embedded in other sites

- System captures both data and processing for reuse

- Community Annotation of/ Interaction Around Problem Specific Virtual Observatory Interfaces

# Conclusion

- Cyber Tools for Large Collaborations need to scale in ways that small project tools do not:
  - Edgeless
  - Evolving
  - Community customizable
  - Cost-effective
- There are design patterns that support such scaling
- Implementation of such tools breaks a 'collaboration + data dissemination' model and can provide significant new value to communities
- Inter-project coordination on interfaces and standards can help guide the CI community

NCSA

# Acknowledgments

NCSA CET Staff & Collaborators
WATERS Network & CI Communities

National Science Foundation
State of Illinois
Office of Naval Research
Department of Energy

## … and Thank You

TRECC

Mid-America Earthquake Center

For more info, see
http://cet.ncsa.uiuc.edu/
http://cet.ncsa.uiuc.edu/publications/

NCSA